

Соловьев С.Ю.

Постановки задач современной информатики

park.glossary.ru/modern/

Задачи статистического анализа

2015 –2022

Статистический анализ – теория конструирования и применения методов выявления причинно-следственных связей.

Математическая статистика – раздел математики, в котором разрабатываются методы регистрации, описания и анализа данных наблюдений с целью получения вероятностно-статистических моделей случайных явлений.

Напоминание: **Задача**

<p>Дано</p> <p><i>Исходные данные</i></p>	<p>Известно</p> <p><i>Свойства исх. данных</i></p>
<p><i>Алгоритм / Метод / Способ / Схема</i></p>	
<p>Требуется</p> <p><i>Результирующие данные</i></p>	<p>такое, что</p> <p><i>Свойства рез. данных</i></p>

VS.

***Математическая
статистика***

Частоты

Выборки

Оценки характеристик

***Теория
вероятностей***

Вероятность

Случайные величины

Характеристики СВ:

математическое ожидание,
дисперсия, мода и др.

функции распределения

Задачи математической статистики

Дано

Выборка ...

Известно

*Теория
вероятностей*

Алгоритм / Метод / Способ / Схема

Требуется

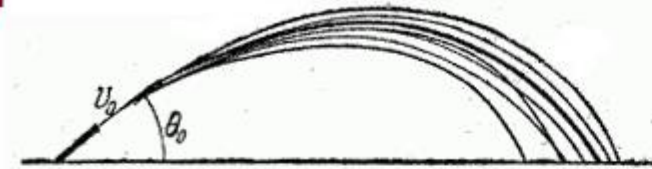
*Оценки
характеристик*

такое, что

*Теория
вероятностей*

Ситуации с неопределенным исходом

Стрельба



Подбрасывание



Стохастические ситуации



**Теория
вероятностей**

Поиск цивилизаций



Основные понятия

Стохастическая ситуация:

- непредсказуемость,
- воспроизводимость,
- устойчивость частот.

Вероятность событий $\longleftarrow \lim$

vs. Аксиоматическое определение вероятности:

$$P : \{ \text{События} \} \rightarrow [0,1] + A1 + A2 + A3 + A4$$

Событие = { Элементарные исходы }

Случайная величина ::=

функция : { Элемент.исходы } \rightarrow { **x** | x - число }

Основные понятия 2

Случайные величины

Дискретные

Распределение
вероятностей

$$\{ (x_k, p_k) \mid k=1, \dots \}$$
$$: p_1 + \dots + p_n = 1$$

Непрерывные

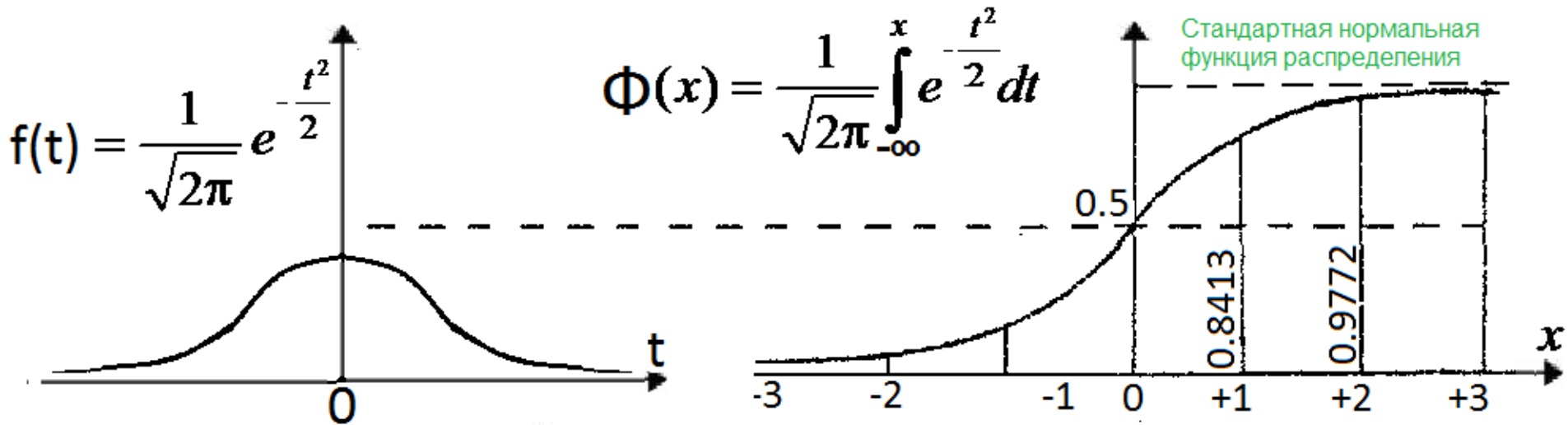
Плотность
вероятностей $f(x) \geq 0$

$$\int_{-\infty}^{+\infty} f(x) dx = 1$$

Функция распределения

$$F(x) = P(X < x) = \begin{cases} \sum_{k: x_k < x} P(X = x_k) & \text{если } X \text{ — дискретная} \\ \int_{-\infty}^x f(t) dt & \text{если } X \text{ — непрерывная} \end{cases}$$

Подзадачи ТВ и МС



Подзадача 1

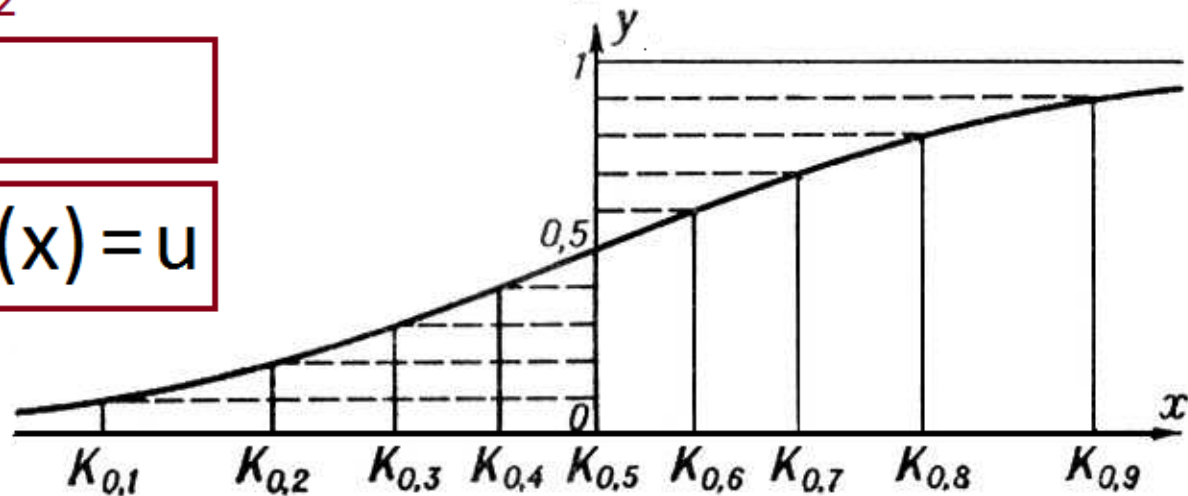
f, x	
--------	--

$\Phi(x)$	
-----------	--

Подзадача 2

f, u	
--------	--

x	$\Phi(x) = u$
-----	---------------



Числовые характеристики СВ X

характеристики “центра”

Математическое ожидание

$$EX = x_1 p_1 + x_2 p_2 + \dots + x_n p_n \quad (\text{диск.})$$

$$EX = \int x f(x) dx \quad // \quad \text{от } -\infty \text{ до } +\infty \quad (\text{непр.})$$

не каждая СВ имеет матем.ожд.

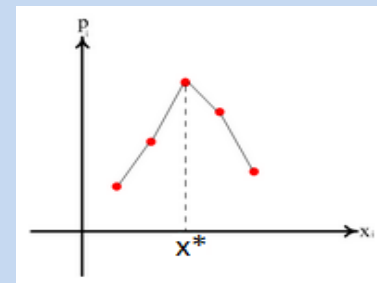
Медиана $\text{med } X ::=$

$$x^* \in \mathbb{R} : P(X \leq x^*) \geq 0.5 \quad \& \quad P(X \geq x^*) \geq 0.5$$

существует, но не единственная

Мода дискретной СВ $\text{mod } X ::=$

$$x^* \in \text{знач. } X : P(X = x_k) \leq P(X = x^*)$$



Числовые характеристики СВ X

характеристики “разброса”

Дисперсия

$$DX = E(X - EX)^2$$

Среднеквадратичное отклонение

$$\sigma = \sqrt{DX} \quad // \quad \sigma^2 = DX$$

Другие

$$E |X - \text{med } X|$$

интерквартильный размах

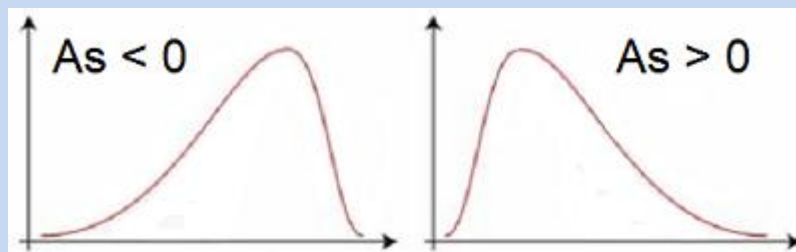
Числовые характеристики СВ X

характеристики “формы”

Коэффициент асимметрии

$$E \left(\frac{X - EX}{\sqrt{DX}} \right)^3$$

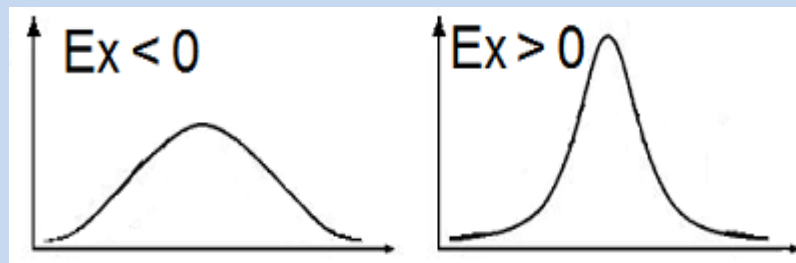
если существует $E|X|^3$



Коэффициент эксцесса

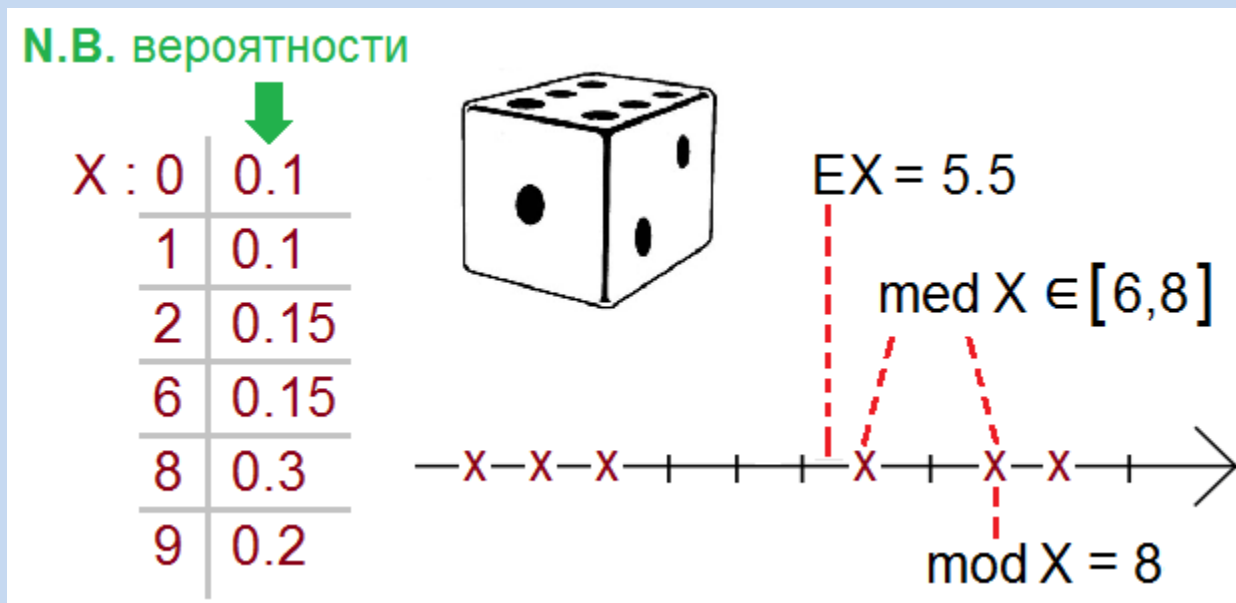
$$E \left(\frac{X - EX}{\sqrt{DX}} \right)^4$$

если существует $E|X|^4$



Числовые характеристики СВ X

пример



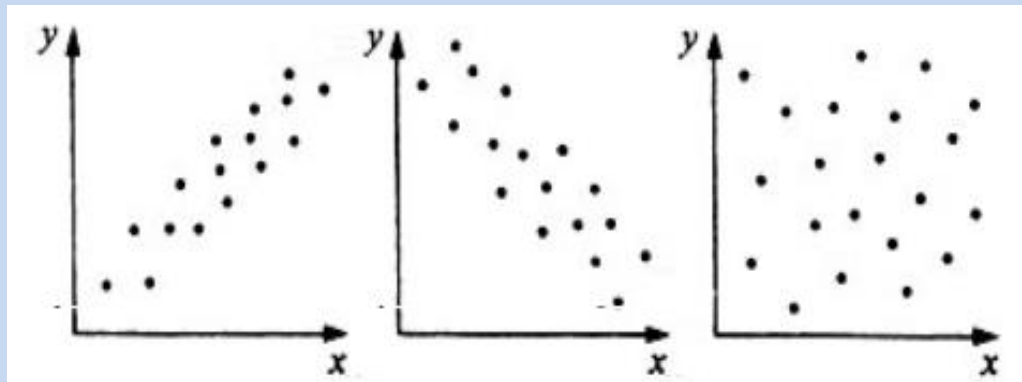
$$DX = 11.25 \quad \sigma = 3.35$$

Числовые характеристики взаимозависимости СВ X и Y

Ковариация $\text{cov}(X, Y) = E(X - EX)(Y - EY)$

Коэффициент корреляции $[-1, +1]$

$$\rho(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{DX DY}}$$



$\rho > 0$

$\rho < 0$

$\rho \approx 0$

Основные понятия математической статистики

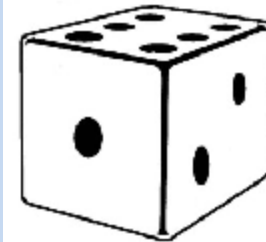
Выборка = X_1, \dots, X_n , где X_i – (а) независимые и
(б) одинаково распределенные СВ.

n – объем выборки

Выборка объемом 14

0, 8, 2, 8, 9, 1, 9, 0, 9, 9, 8, 1, 2, 0

Генеральная совокупность



$X : 0$	0.1
1	0.1
2	0.15
6	0.15
8	0.3
9	0.2

Выборочные числовые характеристики

Выборочное среднее

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

$$\approx EX$$

Выборочная дисперсия

$$S^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

$$\approx DX$$

Выборочное среднеквадратичное отклонение

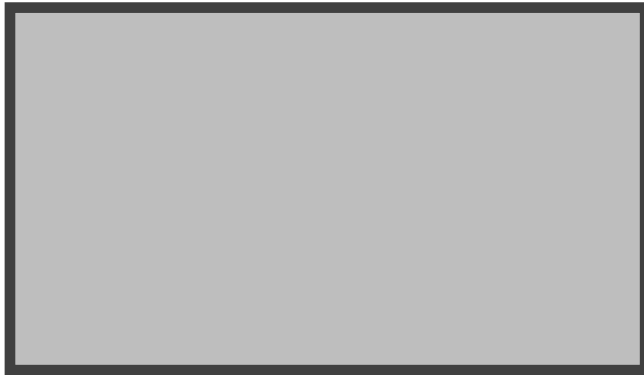
$$S = \sqrt{S^2}$$

$$\approx \sigma$$

Выборочный центральный момент порядка r

$$m_r = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^r$$

Закон больших чисел & Центральная предельная теорема



X_1, \dots, X_n

-- независимые одинаково
распределенные СВ

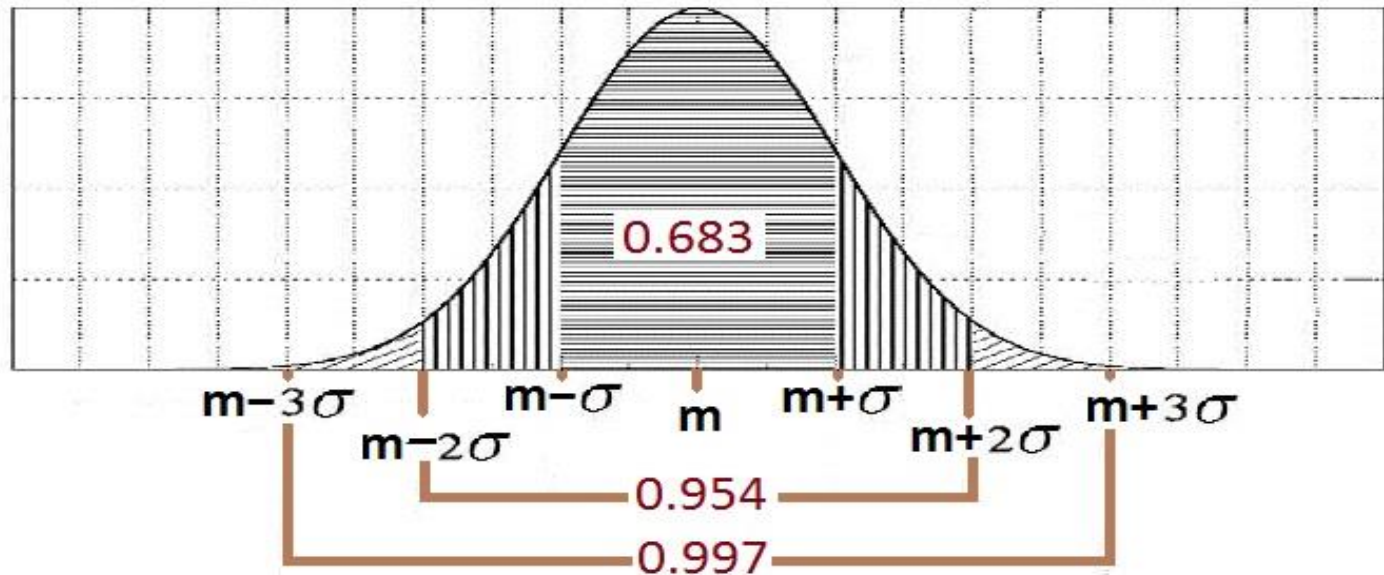
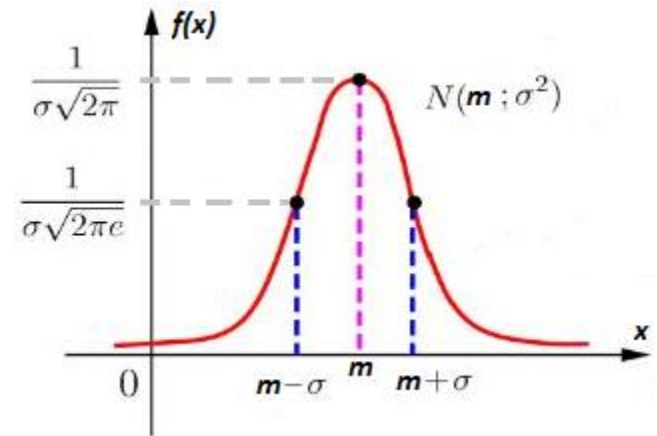
$$P\left(\left|\frac{1}{n}\sum_{k=1}^n X_k - EX\right| < \varepsilon\right) \geq 1 - \frac{DX}{n\varepsilon^2}$$

$$P\left(\sum_{k=1}^n X_k < x\right) \approx \Phi\left(\frac{x - nEX}{\sigma\sqrt{n}}\right)$$

Нормальное распределение $n=1$

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-m)^2}{2\sigma^2}}$$

$m = EX$

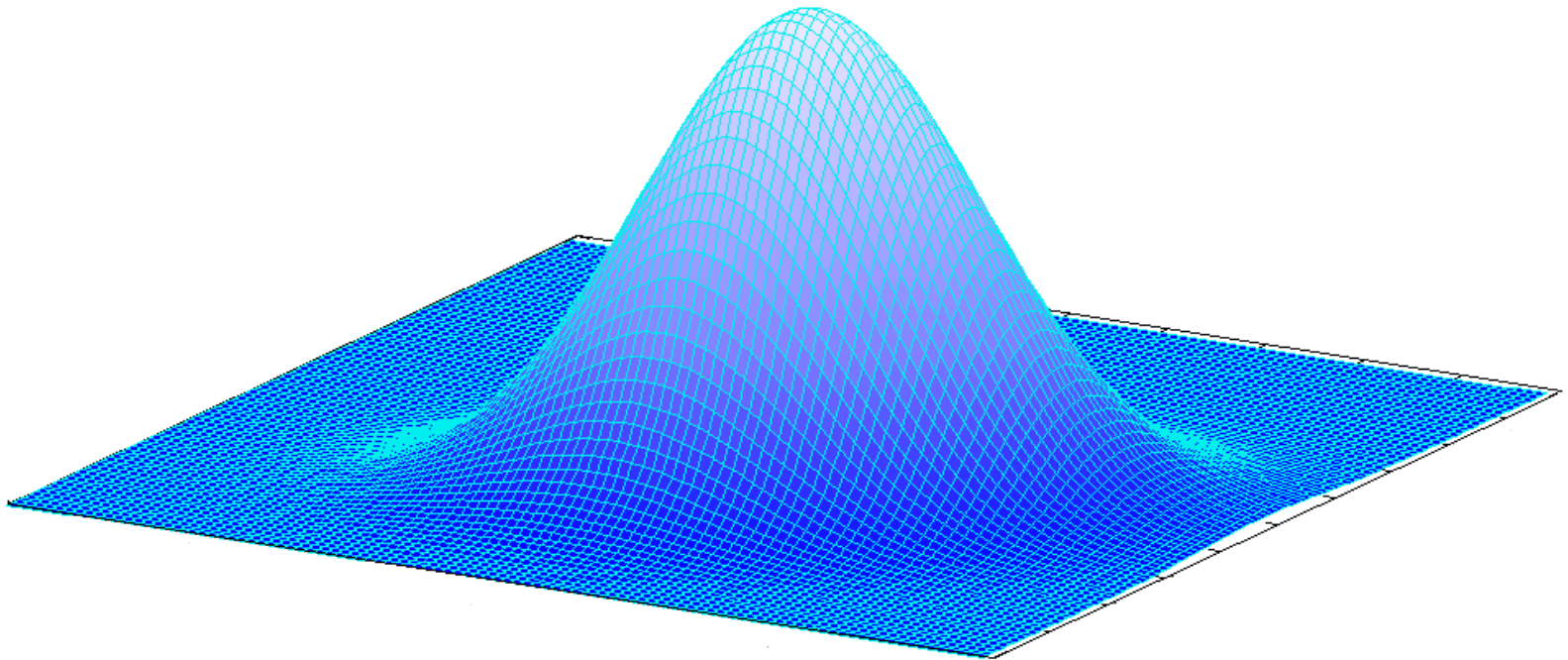


Нормальное распределение $n > 1$

Инфо: многомерная плотность нормального распределения

$$N(x; \mu, \Sigma) = (2^n \pi^n |\Sigma|)^{-0.5} \exp(-0.5(x-\mu)^T \Sigma^{-1}(x-\mu))$$

$\mu \in \mathbb{R}^n$ – матем. ожидание, $\Sigma \in \mathbb{R}^{n \times n}$ – ковар. матрица



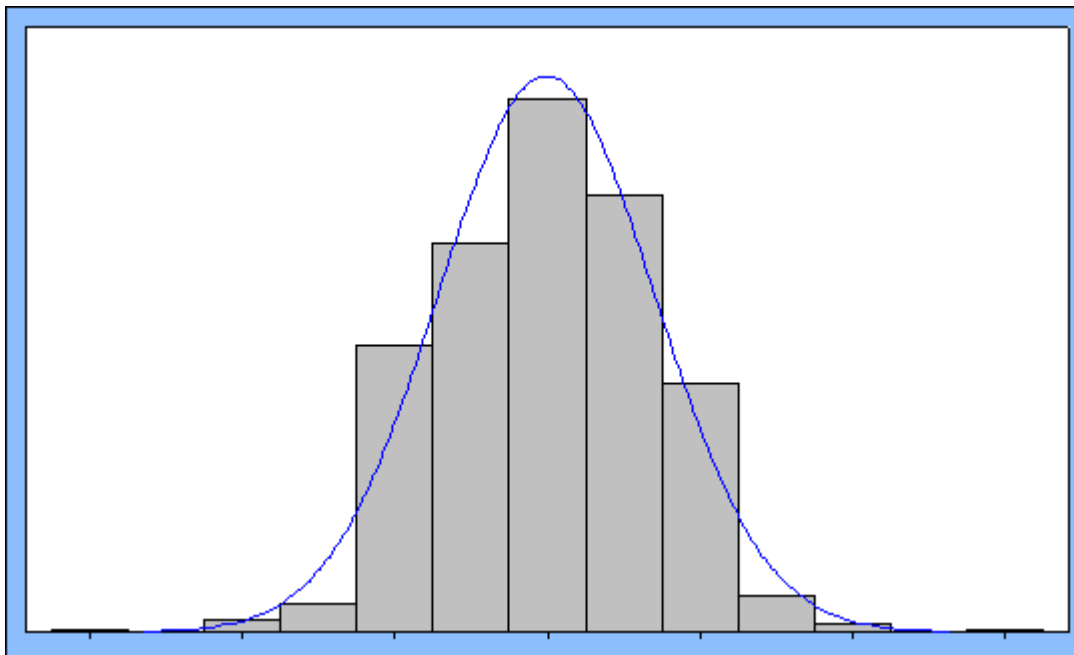
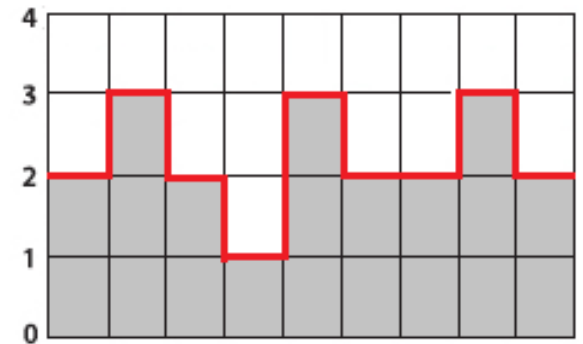
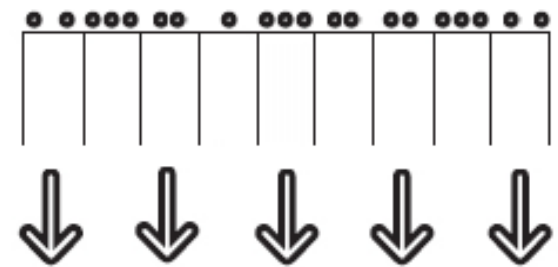
Выборка >> Гистограмма

X_1, \dots, X_{20}
 $n = 20$

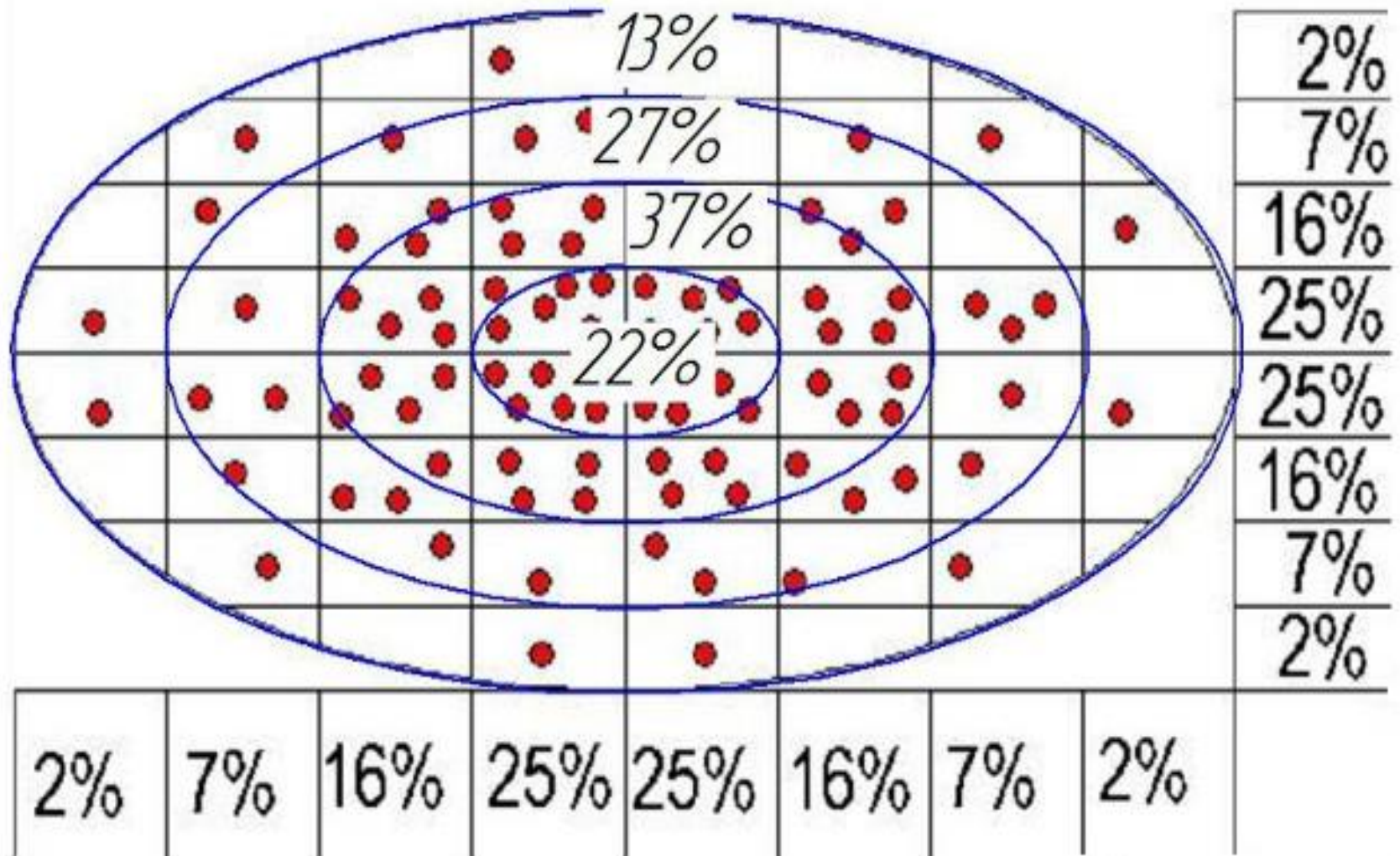


$k = 9$

- $k = [3.3 \lg n + 1]$ – Старджес
- $k = [5 \lg n]$ – Брукс
- $k = [1.72 n^{1/3}]$ – Максимов

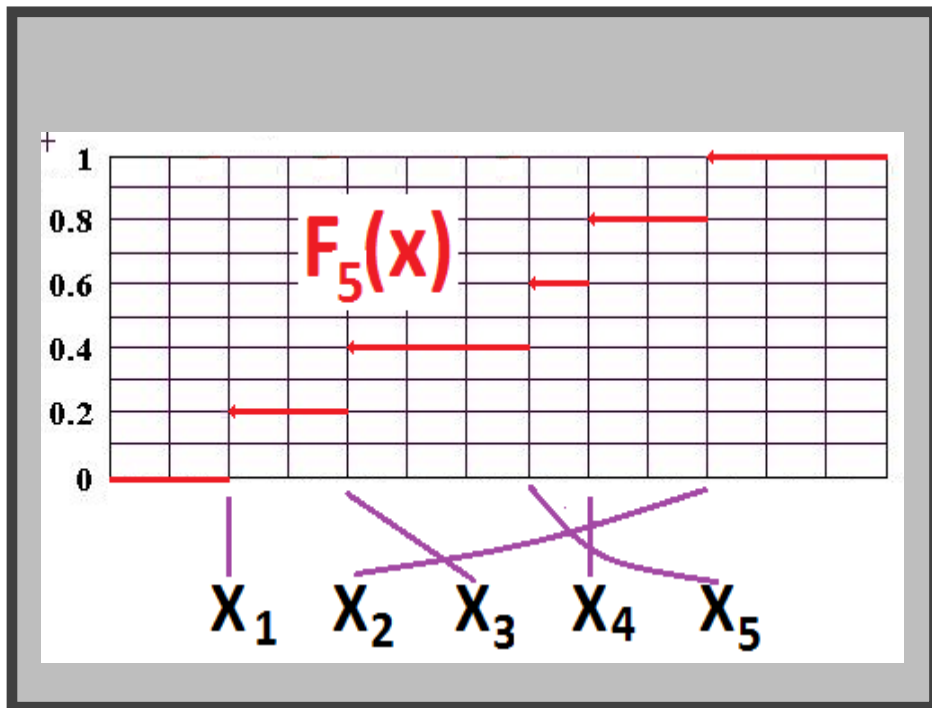


Выборка >> Эллипс рассеяния



[Википедия]

Репрезентативность выборки



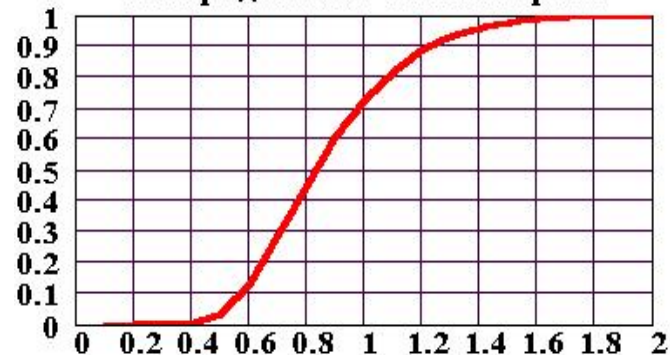
$$P(\max_x |F_n(x) - F(x)| \leq \varepsilon) \geq \gamma$$

$$\Rightarrow K(\gamma) \approx \varepsilon \times \text{sqrt}(n)$$

$$\gamma = 0.95 \quad \varepsilon = 0.035 \quad n \approx 1500$$

$$\gamma = 0.95 \quad \varepsilon = 0.001 \quad n \approx 1849600$$

Распределение Колмогорова



Основные классы задач математической статистики

1. Определение вероятности события.
2. Оценивание числовых характеристик СВ:
 - точечное оценивание;
 - интервальное оценивание;
3. Нахождение закона распределения СВ.
4. Проверка статистических гипотез.
5. Определение регрессионных зависимостей.

Точечные статистические оценки

Точечная статистическая оценка ::=

Функция : Выборка \rightarrow Число,

где **Число** θ – оценка χ характеристики СВ.

Свойства оценок:

- состоятельность: $\theta \rightarrow \chi$ при $n \rightarrow \infty$ по вер.
- несмещаемость: $E\theta = \chi$
- робастность: *устойчивость к выбросам*
- эффективность: *min дисперсии*

Точечные статистические оценки

Выборочное среднее –

(а) состоятельная, (б) **не**смещенная и (в) **не**робастная

оценка математического ожидания СВ.

Выборочная дисперсия –

(а) состоятельная, (б) смещенная и (в) **не**робастная

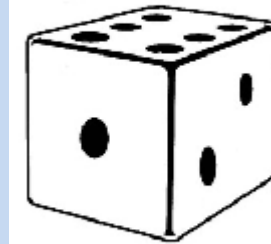
оценка дисперсии СВ.

⇒ *несм.*
$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

Точечные статистические оценки

пример

Объем	500	5000	32000
Частоты	0.11	0.11	0.10
	0.09	0.09	0.10
	0.16	0.16	0.15
	0.12	0.14	0.15
	0.34	0.30	0.30
	0.18	0.20	0.20
	Выб.средн	5.46	5.44
Выб.дисп.	11.54	11.44	11.35



X: 0	0.1
1	0.1
2	0.15
6	0.15
8	0.3
9	0.2

$$EX = 5.5$$

$$DX = 11.25$$

Задача нахождения доверительного интервала

для вероятности события при большом объеме выборки

Дано числа n, p^*, γ	n – объем выборки, $1/(n \times \sqrt{n}) \approx 0$; p – вероятность события A ; (обозн.) p^* – отн. частота события A ; γ – доверительная вероятность;
-----------------------------------	---

Алгоритм*)

числа θ_1 и θ_2	такие, что $P(\theta_1 < p < \theta_2) = \gamma$
-------------------------------	--

- *)
1. Вычислить $\alpha = (1+\gamma)/2$
 2. Найти по таблице число u_α
 3. Вычислить $\varepsilon = u_\alpha \times \sqrt{p^* (1-p^*) / n}$
 4. Положить $\theta_1 = p^* - \varepsilon, \theta_2 = p^* + \varepsilon$

Задача нахождения доверительного интервала

для вероятности события при большом объеме выборки

послесловие

Квантили нормального
распределения $N(0,1)$

$$\Phi(t) = \int_{-\infty}^t \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx$$

$$\Phi(t) = 1 - \Phi(-t)$$

$$\Phi(u_\alpha) = \alpha$$

?	u_α	?	u_α	?	u_α
.50	0	.91	1.341	.995	2.576
.55	.126	.92	1.405	.999	3.090
.60	.253	.93	1.476	.9995	3.291
.65	.385	.94	1.555	.9999	3.719
.70	.524	.95	1.645	.99995	3.891
.75	.674	.96	1.751	.99999	4.265
.80	.842	.97	1.881	.999995	4.417
.85	1.036	.98	2.054	.999999	4.753
.90	1.282	.99	2.326	.9999999	5.199

$$n = 1600, p^* = 0.2, \gamma = 0.95$$

$$0.18 < p < 0.22 \text{ с надежностью } 0.95$$

Задача нахождения доверительного интервала

для математического ожидания при большом объеме выборки

Дано числа n , γ выборка X_1, \dots, X_n	n – объем выборки, $n > 30$; γ – доверительная вероятность; EX – математическое ожидание СВ X ; (обозн.)
--	--

Алгоритм*)

числа θ_1 и θ_2	такие, что $P(\theta_1 < EX < \theta_2) = \gamma$
-------------------------------	---

- *)
1. Вычислить $\alpha = (1+\gamma)/2$ и найти по таблице число u_α
 2. Вычислить выборочные характеристики \bar{X} и S ;
 3. Вычислить $\varepsilon = S \times u_\alpha / \text{sqrt}(n)$
 4. Положить $\theta_1 = \bar{X} - \varepsilon$, $\theta_2 = \bar{X} + \varepsilon$

Задача нахождения доверительного интервала

для среднеквадратичного отклонения

при большом объеме выборки

Дано числа n , γ выборка X_1, \dots, X_n	n – объем выборки; γ – доверительная вероятность; DX – среднеквадратичное отклонение СВ X ;
--	--

Алгоритм*)

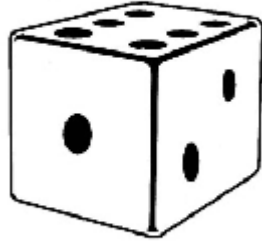
числа θ_1 и θ_2	такие, что $P(\theta_1 < DX < \theta_2) = \gamma$
-------------------------------	---

- *)
1. Вычислить $\alpha = (1+\gamma)/2$ и найти по таблице число u_α
 2. Вычислить выборочные характеристики m_4 и S ;
 3. Вычислить $\varepsilon = u_\alpha \times \sqrt{m_4 / S^4 - 1} / \sqrt{4 \times n}$
 4. Положить $\theta_1 = S \times (1 - \varepsilon)$, $\theta_2 = S \times (1 + \varepsilon)$

Интервальные статистические оценки

пример ($\gamma = 0.95$)

Объем	500	32000
Частоты	0.11 ± 0.027	0.10 ± 0.003
	0.09 ± 0.026	0.10 ± 0.003
	0.16 ± 0.032	0.15 ± 0.004
	0.12 ± 0.028	0.15 ± 0.004
	0.34 ± 0.042	0.30 ± 0.005
	0.18 ± 0.038	0.20 ± 0.004
	0.18 ± 0.038	0.20 ± 0.004
Выб.средн	5.46 ± 0.30	5.46 ± 0.037
Выб.скоткл	3.40 ± 0.11	3.36 ± 0.013



X: 0	0.1
1	0.1
2	0.15
6	0.15
8	0.3
9	0.2

$$EX = 5.5$$

$$DX = 11.25$$

$$\sigma = 3.35$$

Проверка статистических гипотез

Статистическая гипотеза – предположение о свойствах генерального распределения, проверяемые статистическими методами.

Критерий значимости – правило проверки статистической гипотезы.

Статистика критерия – функция $Z : \text{Выборка} \rightarrow \text{Число}$, по значениям которой судят о справедливости статистической гипотезы.

Уровень значимости α :

событие с вер. $\alpha \equiv$ событие практически невозможное,
событие с вер. $1-\alpha \equiv$ событие практически достоверное.

H_0 – проверяемая гипотеза **vs.** H_1 – альтернативная гипотеза.

Критическая область критерия $V_k : P(Z \in V_k | H_0) = \alpha$

Область допустимых значений = $V \setminus V_k$

≈Схема проверки статистических гипотез

Этап 1. Сформулировать гипотезы H_0 и H_1 .

Этап 2. Выбрать уровень значимости α .

Этап 3. Выбрать статистику $Z = Z(\text{Выборка})$ для проверки H_0 .

Этап 4. Найти распределение $F(z | H_0)$: хи-квадрат, Стьюдента и др.

Этап 5. Построить критическую область V_k : $P(H_1 | H_0) \oplus P(H_0 | H_1) \leq \alpha$

Этап 6. Вычислить значение $Z_B = Z(\text{Выборка})$.

Этап 7. На уровне доверия $1-\alpha$:

если $Z_B \in V_k$, то H_0 отклоняется.

если $Z_B \in V \setminus V_k$, то H_0 принимается.



Вычисление статистики Z

$\chi^2(n)$ Распределение хи-квадрат с n степенями свободы
– распределение СВ $X_1^2 + \dots + X_n^2$, где $X_i \sim N(0;1)$.

$S(n)$ Распределение Стьюдента с n степенями свободы
– распределение СВ $X / \sqrt{Y / n}$, где $X \sim N(0;1)$, $Y \sim \chi^2(n)$

$F(n;m)$ Распределение Фишера с n и m степенями свободы
– распределение СВ $(X / n) / (Y / m)$, где $X \sim \chi^2(n)$, $Y \sim \chi^2(m)$

Таблица $\chi^2(n)$ | Таблица $t(n)$ | Таблица $F(n;m)$

Проверка гипотезы о равенстве дисперсий нормально распределенных ген.совокупностей

Выборки

X_1, \dots, X_n

Y_1, \dots, Y_m

число α

Обе выборки получены независимо

из нормальных ген. совокупностей X или Y ;

$H_0 : DX = DY$; $H_1 : DX \neq DY$;

α -- уровень значимости

Критерий*)

либо H_0 принимается, либо H_0 отклоняется

- *) 1 . Вычислить выборочные несмещаемые дисп. S_X^2 и S_Y^2
- 2 . Вычислить статистику $F_B = S_X^2 / S_Y^2$;
- 3 . Найти по таблице значение $K = F_{1-\alpha/2}(m-1;n-1)$;
- 4 . Если $F_B < K$, то H_0 принимается, иначе H_0 отклоняется

Проверка гипотезы о равенстве вероятностей двух событий при больших объемах выборок

Отн. частоты

p_A^*, p_B^*

число α

p_A^*, p_B^* – относительные частоты событий **A** и **B** ;; получены из выборок большого объема.

p_A, p_B – вероятности событий **A** и **B** ; ; ;

$H_0 : p_A = p_B$; ; ; ; ; $H_1 : p_A \neq p_B$; ; ; ; ; ; ; ; ;

α -- уровень значимости

Упрощенный критерий*) .

Используется доверительный интервал.

либо H_0 принимается, либо H_0 отклоняется

*) 1 . Вычислить доверительный интервал – числа θ_1 и θ_2 :

$$P(\theta_1 < p_A < \theta_2) = 1 - \alpha$$

2 . Если $\theta_1 < p_B^* < \theta_2$,

то H_0 принимается, иначе H_0 отклоняется

Проверка гипотезы о законе распределения генеральной совокупности

Выборка	X_1, \dots, X_n	$F(x; \theta_1, \dots, \theta_k)$ – закон распределения ;
Числа	$\theta_1, \dots, \theta_k$	H_0 : распределение ген. совокупности X есть $F(x; \theta_1, \dots, \theta_k)$; ; ; ; ; ; ; ; ; ;
Число	α	α – уровень значимости

Критерий хи-квадрат Пирсона*) (идея)

либо H_0 принимается, либо H_0 отклоняется

*) 1. Выборка \Rightarrow Последовательность интервалов $\Delta_1, \dots, \Delta_m$

2. Выч. n_i – к-во X_j в Δ_i и $p_i = P(x \in \Delta_i)$

3. Вычислить статистику

$$Z = \sum_{i=1}^m \frac{(n_i - np_i)^2}{np_i}$$

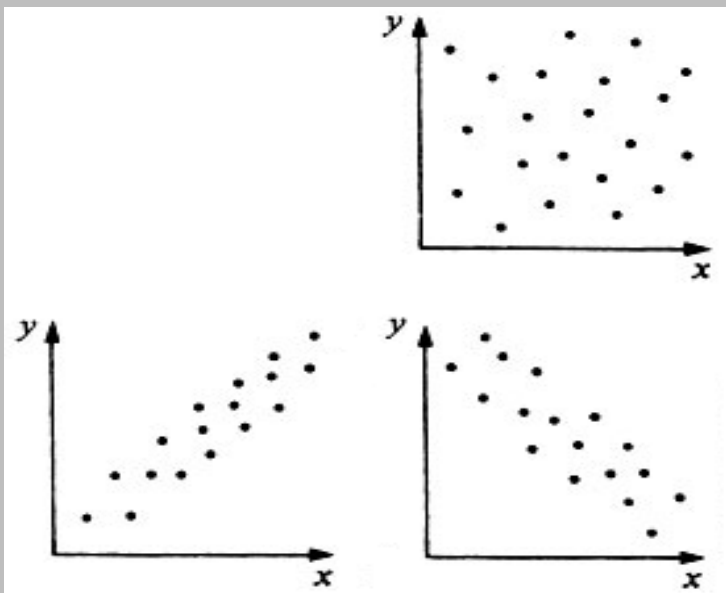
4. Если $0 < Z < \chi^2_{1-\alpha}(m-k)$,

то H_0 принимается, иначе отклоняется

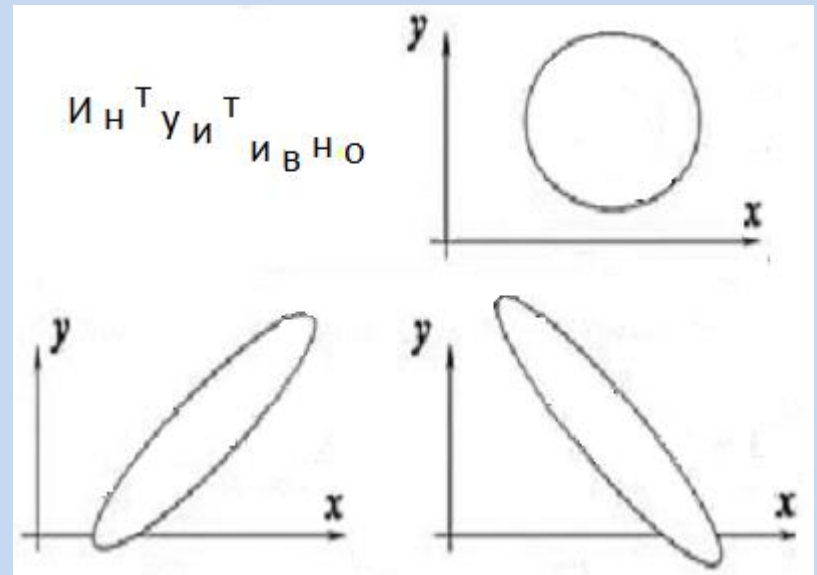
Корреляционный анализ

-- раздел математической статистики, исследующий свойства статистических зависимостей между СВ.

Двумерная выборка
 $(X_1, Y_1), \dots, (X_n, Y_n)$



Генеральная совокупность =
двумерная СВ: $\{ (x_i, y_i) \}$
Распределение: $\{ ((x_i, y_i), p_i) \}$



Порождает СВ X и Y , и вопрос о \exists зависимости между X и Y .

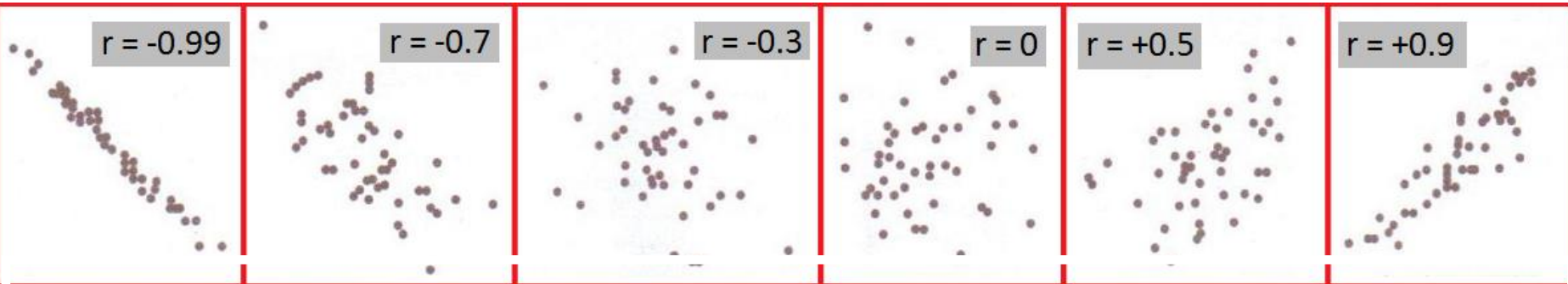
Коэффициент корреляции

$(X_1, Y_1), \dots, (X_n, Y_n) \Rightarrow$ Выборочный коэффициент корреляции:

$$r_{XY} = \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{S_X S_Y}$$

Коэффициент корреляции двух СВ:

$$\rho(X, Y) = \frac{E(X - EX)(Y - EY)}{\sqrt{DX DY}}$$



r_{XY} vs. $\rho(X, Y)$

Грубая проверка гипотезы о независимости случайных величин

Если $|r_{XY}| \times \sqrt{n-1} > 2.5,$

то X, Y зависимы с уровнем значимости 0.005

Если $|r_{XY}| \times \sqrt{n-1} > 3.0,$

то X, Y зависимы с уровнем значимости 0.001

Регрессионный анализ

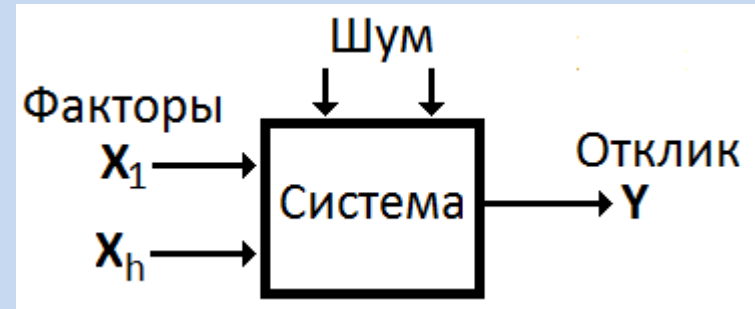
-- раздел математической статистики, исследующий статистические зависимости между двумя и более СВ.

Двумерная выборка
 $(X_1, Y_1), \dots, (X_n, Y_n)$



$$y(x) = B_0 + B_1 x$$

→ эмпирическая простая линейная регрессия



Зависимость Y от X_1, \dots, X_h ?

1. $h = 1$
2. Зависимость EY от X , т.е. цель -- поиск функции

$$y(x) = M_x Y$$

3. $y(x) = \beta_0 + \beta_1 x$

→ теоретическая простая линейная регрессия

Построение эмпирической простой линейной регрессии

Метод наименьших квадратов:

Дано $(X_1, Y_1), \dots, (X_n, Y_n)$. Требуется B_1 и B_2 такие, что

$$\sum_{i=1..n} (Y_i - B_0 - B_1 x_i)^2 \rightarrow \min$$

Решение: $B_1 = r_{XY} \frac{S_Y}{S_X}, \quad B_0 = \bar{Y} - \bar{X} B_1$

+ Условия Гаусса – Маркова (1777-1855; 1856-1922)

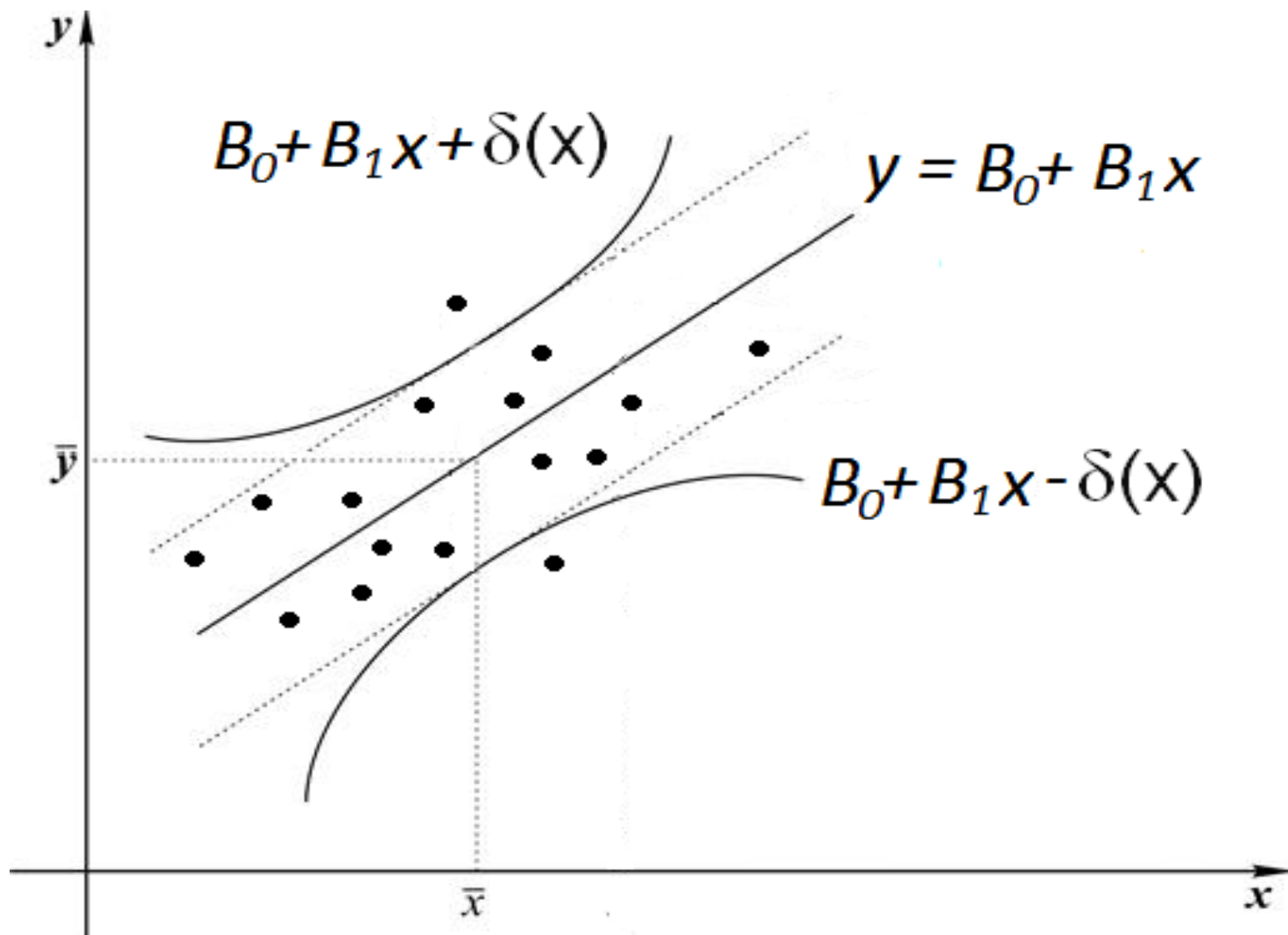
Отклонения $Y_i - B_0 - B_1 X_i$

(а) нормально распределены,

(б) независимы,

(в) имеют нулевые м.ожидания и равные дисперсии.

Коридор ошибок



В о п р о с ы?

soloviev@glossary.ru

Соловьев С.Ю. Постановки задач современной информатики.
www.park.glossary.ru